

---

# CHAPTER 7

---

## Speech Coding

*T*ransmission of speech using data networks requires the speech signal to be digitally encoded. Voice over IP has become very popular because of the Internet, where bandwidth limitations make it necessary to compress the speech signal. Digital storage of audio signals, which can result in higher quality and smaller size than the analog counterpart, is commonplace in compact discs, digital video discs, and MP3 files. Many spoken language systems also use coded speech for efficient communication. For these reasons we devote a chapter to speech and audio coding techniques.

Rather than exhaustively cover all the existing speech and audio coding algorithms we uncover their underlying technology and enumerate some of the most popular standards. The coding technology discussed in this chapter has a strong link to both speech recognition and speech synthesis. For example, the speech synthesis algorithms described in Chapter 16 use many techniques described here.

## 7.1. SPEECH CODERS ATTRIBUTES

How do we compare different speech or audio coders? We can refer to a number of factors, such as signal bandwidth, bit rate, quality of reconstructed speech, noise robustness, computational complexity, delay, channel-error sensitivity and standards.

Speech signals can be bandlimited to 10 kHz without significantly affecting the hearer's perception. The telephone network limits the bandwidth of speech signals to between 300 and 3400 Hz, which gives *telephone speech* a lower quality. Telephone speech is typically sampled at 8 kHz. The term *wideband speech* is used for a bandwidth of 50–7000 Hz and a sampling rate of 16 kHz. Finally, *audio coding* is used in dealing with high-fidelity audio signals, in which case the signal is sampled at 44.1 kHz.

Reduction in bit rate is the primary purpose of speech coding. The previous bit stream can be compressed to a lower rate by removing redundancy in the signal, resulting in savings in storage and transmission bandwidth. If only redundancy is removed, the original signal can be recovered exactly (*lossless* compression). In *lossy* compression, the signal cannot be recovered exactly, though hopefully it will sound similar to the original.

Depending on system and design constraints, fixed-rate or variable-rate speech coders can be used. Variable-rate coders are used for non-real time applications, such as voice storage (silence can be coded with fewer bits than fricatives, which in turn use fewer bits than vowels), or for packet voice transmissions, such as CDMA cellular for better channel utilization. Transmission of coded speech through a noisy channel may require devoting more bits to channel coding and fewer to source coding. For most real-time communication systems, a maximum bit rate is specified.

The quality of the reconstructed speech signal is a fundamental attribute of a speech coder. Bit rate and quality are intimately related: the lower the bit rate, the lower the quality. While the bit rate is inherently a number, it is difficult to quantify the quality. The most widely used measure of quality is the *Mean Opinion Score* (MOS) [25], which is the result of averaging opinion scores for a set of between 20 and 60 untrained subjects. Each listener characterizes each set of utterances with a score on a scale from 1 (unacceptable quality) to 5 (excellent quality), as shown in Table 7.1. An MOS of 4.0 or higher defines *good* or *toll* quality, where the reconstructed speech signal is generally indistinguishable from the original signal. An MOS between 3.5 and 4.0 defines *communication* quality, which is sufficient for telephone communications. We show in Section 7.2.1 that if each sample is quantized with 16 bits, the resulting signal has *toll* quality (essentially indistinguishable from the unquantized signal). See Chapter 16 for more details on perceptual quality measurements.

**Table 7.1** Mean Opinion Score (MOS) is a numeric value computed as an average for a number of subjects, where each number maps to the above subjective quality.

Excellent	Good	Fair	Poor	Bad
5	4	3	2	1

Another measure of quality is the *signal-to-noise ratio* (SNR), defined as the ratio between the signal's energy and the noise's energy in terms of dB:

$$SNR = \frac{\sigma_x^2}{\sigma_e^2} = \frac{E\{x^2[n]\}}{E\{e^2[n]\}} \quad (7.1)$$

The MOS rating of a codec on noise-free speech is often higher than its MOS rating for noisy speech. This is generally caused by specific assumptions in the speech coder that tend to be violated when a significant amount of noise is present in the signal. This phenomenon is more accentuated for lower-bit-rate coders that need to make more assumptions.

The computational complexity and memory requirements of a speech coder determine the cost and power consumption of the hardware on which it is implemented. In most cases, real-time operation is required at least for the decoder. Speech coders can be implemented in *inexpensive Digital Signal Processors* (DSP) that form part of many consumer devices, such as answering machines and DVD players, for which storage tends to be relatively more expensive than processing power. DSPs are also used in cellular phones because bit rates are limited.

All speech coders have some delay, which, if excessive, can affect the dynamics of a two-way communication. For instance, delays over 150 ms can be unacceptable for highly interactive conversations. Coder delay is the sum of different types of delay. The first is the *algorithmic delay* arising because speech coders usually operate on a block of samples, called a *frame*, which needs to be accumulated before processing can begin. Often the speech coder requires some additional *look-ahead* beyond the frame to be encoded. The *computational delay* is the time that the speech coder takes to process the frame. For real-time operation, the computational delay has to be smaller than the algorithmic delay. A block of bits is generally assembled by the encoder prior to transmission, possibly to add error-correction properties to the bit stream, which cause *multiplexing delay*. Finally, there is the *transmission delay*, due to the time it takes for the frame to traverse the channel. The decoder will incur a *decoder delay* to reconstruct the signal. In practice, the total delay of many speech coders is at least three frames.

If the coded speech needs to be transmitted over a channel, we need to consider possible channel errors, and our speech decoder should be insensitive to at least some of them. There are two types of errors: random errors and burst errors, and they could be handled differently. One possibility to increase the robustness against such errors is to use channel coding techniques, such as those proposed in Chapter 3. Joint source and channel coding allows us to find the right combination of bits to devote to speech coding with the right amount devoted to channel coding, adjusting this ratio adaptively depending on the channel. Since channel coding will only reduce the number of errors, and not eliminate them, graceful degradation of speech quality under channel errors is typically a design factor for speech coders. When the channel is the Internet, complete frames may be missing because they have not arrived in time. Therefore, we need techniques that degrade gracefully with missing frames.

## 7.2. SCALAR WAVEFORM CODERS

In this section we describe several waveform coding techniques, such as linear PCM,  $\mu$ -law, and A-law PCM, APCM, DPCM, DM, and ADPCM, that quantize each sample using scalar quantization. These techniques attempt to approximate the waveform, and, if a large enough bit rate is available, will get arbitrarily close to it.

### 7.2.1. Linear Pulse Code Modulation (PCM)

Analog-to-digital converters perform both sampling and quantization simultaneously. To better understand how this process affects the signal it's better to study them separately. We analyzed the effects of sampling in Chapter 5, so now we analyze the effects of quantization, which encodes each sample with a fixed number of bits. With  $B$  bits, it is possible to represent  $2^B$  separate quantization levels. The output of the quantizer  $\hat{x}[n]$  is given by

$$\hat{x}[n] = Q\{x[n]\} \quad (7.2)$$

Linear *Pulse Code Modulation* (PCM) is based on the assumption that the input discrete signal  $x[n]$  is bounded

$$|x[n]| \leq X_{\max} \quad (7.3)$$

and that we use *uniform quantization* with quantization step size  $\Delta$  which is constant for all levels  $x_i$

$$x_i - x_{i-1} = \Delta \quad (7.4)$$

The input/output characteristics are shown by Figure 7.1 for the case of a 3-bit uniform quantizer. The so-called *mid-riser* quantizer has the same number of positive and negative levels, whereas the *mid-tread* quantizer has one more negative than positive levels. The code  $c[n]$  is expressed in two's complement representation, which for Figure 7.1 varies between  $-4$  and  $+3$ . For the mid-riser quantizer the output  $\hat{x}[n]$  can be obtained from the code  $c[n]$  through

$$\hat{x}[n] = \text{sign}(c[n]) \frac{\Delta}{2} + c[n]\Delta \quad (7.5)$$

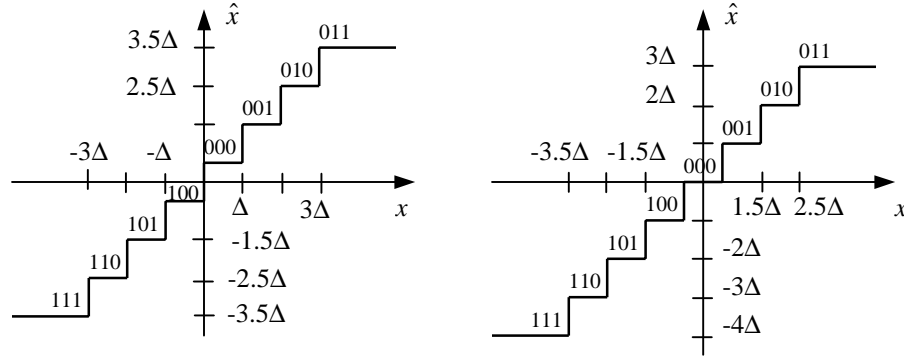
and for the mid-tread quantizer

$$\hat{x}[n] = c[n]\Delta \quad (7.6)$$

which is often used in computer systems that use two's complement representation.

There are two independent parameters for a uniform quantizer: the number of levels  $N = 2^B$ , and the step size  $\Delta$ . Assuming Eq. (7.3), we have the relationship

$$2X_{\max} = \Delta 2^B \quad (7.7)$$



**Figure 7.1** Three-bit uniform quantization characteristics: (a) mid-riser, (b) mid-tread.

In quantization, it is useful to express the relationship between the unquantized sample  $x[n]$  and the quantized sample  $\hat{x}[n]$  as

$$\hat{x}[n] = x[n] + e[n] \quad (7.8)$$

with  $e[n]$  being the quantization noise. If we choose  $\Delta$  and  $B$  to satisfy Eq. (7.7), then

$$-\frac{\Delta}{2} \leq e[n] \leq \frac{\Delta}{2} \quad (7.9)$$

While there is obviously a deterministic relationship between  $e[n]$  and  $x[n]$ , it is convenient to assume a probabilistic model for the quantization noise:

1.  $e[n]$  is white:  $E\{e[n]e[n+m]\} = \sigma_e^2 \delta[m]$
2.  $e[n]$  and  $x[n]$  are uncorrelated:  $E\{x[n]e[n+m]\} = 0$
3.  $e[n]$  is uniformly distributed in the interval  $(-\Delta/2, \Delta/2)$

These assumptions are unrealistic for some signals, except in the case of speech signals, which rapidly fluctuate between different quantization levels. The assumptions are reasonable if the step size  $\Delta$  is a small enough, or alternatively the number of levels is large enough (say more than  $2^6$ ).

The variance of such uniform distribution (see Chapter 3) is

$$\sigma_e^2 = \frac{\Delta^2}{12} = \frac{X_{\max}^2}{3 \times 2^{2B}} \quad (7.10)$$

after using Eq. (7.7). The SNR is given by

$$SNR(dB) = 10 \log_{10} \left( \frac{\sigma_x^2}{\sigma_e^2} \right) = (20 \log_{10} 2)B + 10 \log_{10} 3 - 20 \log_{10} \left( \frac{X_{\max}}{\sigma_x} \right) \quad (7.11)$$

which implies that each bit contributes to 6 dB of SNR, since  $20 \log_{10} 2 \cong 6$ .

Speech samples can be approximately described as following a *Laplacian* distribution [40]

$$p(x) = \frac{1}{\sqrt{2}\sigma_x} e^{-\frac{\sqrt{2}|x|}{\sigma_x}} \quad (7.12)$$

and the probability of  $x$  falling outside the range  $(-4\sigma_x^2, 4\sigma_x^2)$  is 0.35%. Thus, using  $X_{\max} = 4\sigma_x$ ,  $B = 7$  bits in Eq. (7.11) results in an *SNR* of 35 dB, which would be acceptable in a communications system. Unfortunately, signal energy can vary over 40 dB, due to variability from speaker to speaker as well as variability in transmission channels. Thus, in practice, it is generally accepted that 11 bits are needed to achieve an *SNR* of 35dB while keeping the clipping to a minimum.

Digital audio stored in computers (Windows WAV, Apple AIF, Sun AU, and SND formats among others) use 16-bit linear PCM as their main format. The *Compact Disc-Digital Audio* (CD-DA or simply CD) also uses 16-bit linear PCM. Invented in the late 1960s by James T. Russell, it was launched commercially in 1982 and has become one of the most successful examples of consumer electronics technology: there were about 700 million audio CD players in 1997. A CD can store up to 74 minutes of music, so the total amount of digital data that must be stored on a CD is 44,100 samples/(channel\*second) \* 2 bytes/sample \* 2 channels \* 60 seconds/minute \* 74 minutes = 783,216,000 bytes. This 747 MB are stored in a disk only 12 centimeters in diameter and 1.2 mm thick. CD-ROMs can record only 650 MB of computer data because they use the remaining bits for error correction.

### 7.2.2. $\mu$ -law and A-law PCM

Human perception is affected by *SNR*, because adding noise to a signal is not as noticeable if the signal energy is large enough. Ideally, we want *SNR* to be constant for all quantization levels, which requires the step size to be proportional to the signal value. This can be done by using a logarithmic *comparator*<sup>1</sup>

$$y[n] = \ln|x[n]| \quad (7.13)$$

followed by a uniform quantizer on  $y[n]$  so that

$$\hat{y}[n] = y[n] + \varepsilon[n] \quad (7.14)$$

and, thus,

$$\hat{x}[n] = \exp\{\hat{y}[n]\}\text{sign}\{x[n]\} = x[n]\exp\{\varepsilon[n]\} \quad (7.15)$$

after using Eq. (7.13) and (7.14). If  $\varepsilon[n]$  is small, then Eq. (7.15) can be expressed as

$$\hat{x}[n] \cong x[n](1 + \varepsilon[n]) = x[n] + x[n]\varepsilon[n] \quad (7.16)$$

<sup>1</sup> A *comparator* is a nonlinear function that compands one part of the  $x$ -axis.

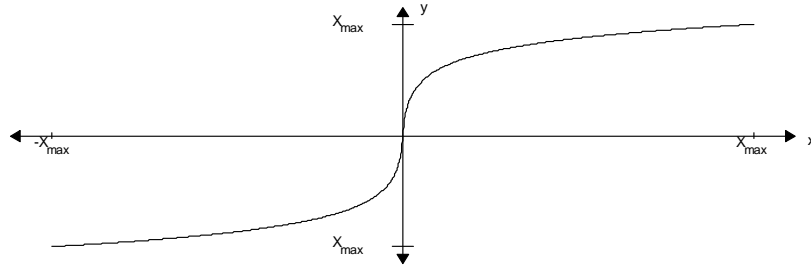
and, thus, the  $SNR = 1/\sigma_e^2$  is constant for all levels. This type of quantization is not practical, because an infinite number of quantization steps would be required. An approximation is the so-called  $\mu$ -law [51]:

$$y[n] = X_{\max} \frac{\log \left[ 1 + \mu \frac{|x[n]|}{X_{\max}} \right]}{\log[1 + \mu]} \text{sign}\{x[n]\} \quad (7.17)$$

which is approximately logarithmic for large values of  $x[n]$  and approximately linear for small values of  $x[n]$ . A related compander called A-law is also used

$$y[n] = X_{\max} \frac{1 + \log \left[ \frac{A|x[n]|}{X_{\max}} \right]}{1 + \log A} \text{sign}\{x[n]\} \quad (7.18)$$

which has greater resolution than  $\mu$ -law for small sample values, but a range equivalent to 12 bits. In practice, they both offer similar quality. The  $\mu$ -law curve can be seen in Figure 7.2.



**Figure 7.2** Nonlinearity used in the  $\mu$ -law compression.

In 1972 the ITU-T<sup>2</sup> recommendation G.711 standardized telephone speech coding at 64 kbps for digital transmission of speech through telephone networks. It uses 8 bits per sample and an 8-kHz sampling rate with either  $\mu$ -law or A-law. In North America and Japan,  $\mu$ -law with  $\mu = 255$  is used, whereas, in the rest of the world, A-law with  $A = 87.56$  is used. Both compression characteristics are very similar and result in an approximate SNR of 35 dB. Without the logarithmic compressor, a uniform quantizer requires approximately 12 bits per sample to achieve the same level of quality. All the speech coders for telephone speech described in this chapter use G.711 as a baseline reference, whose quality is considered *toll*,

<sup>2</sup> The International Telecommunication Union (ITU) is a part of the United Nations Economic, Scientific and Cultural Organization (UNESCO). ITU-T is the organization within ITU responsible for setting global telecommunication standards. Within ITU-T, Study Group 15 (SG15) is responsible for formulating speech coding standards. Prior to 1993, telecommunication standards were set by the *Comité Consultatif International Téléphonique et Télégraphique* (CCITT), which was reorganized into the ITU-T that year.

and an MOS of about 4.0. G.711 is used by most digital central office switches, so that when you make a telephone call using your plain old telephone service (POTS), your call is encoded with G.711. G.711 has an MOS of about 4.3.

### 7.2.3. Adaptive PCM

When quantizing speech signals we confront a dilemma. On the one hand, we want the quantization step size to be large enough to accommodate the maximum peak-to-peak range of the signal and avoid clipping. On the other hand, we need to make the step size small to minimize the quantization noise. One possible solution is to adapt the step size to the level of the input signal.

The basic idea behind *Adaptive PCM* (APCM) is to let the step size  $\Delta[n]$  be proportional to the standard deviation of the signal  $\sigma[n]$ :

$$\Delta[n] = \Delta_0 \sigma[n] \quad (7.19)$$

An equivalent method is to use a fixed quantizer but have a time-varying gain  $G[n]$ , which is inversely proportional to the signal's standard deviation

$$G[n] = G_0 / \sigma[n] \quad (7.20)$$

Estimation of the signal's variance, or short-time energy, is typically done by low-pass filtering  $x^2[n]$ . With a first-order IIR filter, the variance  $\sigma^2[n]$  is computed as

$$\sigma^2[n] = \alpha \sigma^2[n-1] + (1-\alpha)x^2[n-1] \quad (7.21)$$

with  $\alpha$  controlling the time constant of the filter  $T = -1 / (F_s \ln \alpha)$ ,  $F_s$  the sampling rate, and  $0 < \alpha < 1$ . In practice,  $\alpha$  is chosen so that the time constant ranges between 1 ms ( $\alpha = 0.88$  at 8 kHz) and 10 ms ( $\alpha = 0.987$  at 8 kHz).

Alternatively,  $\sigma^2[n]$  can be estimated from the past  $M$  samples:

$$\sigma^2[n] = \frac{1}{M} \sum_{m=n-M}^{n-1} x^2[m] \quad (7.22)$$

In practice, it is advantageous to set limits on the range of values of  $\Delta[n]$  and  $G[n]$ :

$$\Delta_{\min} \leq \Delta[n] \leq \Delta_{\max} \quad (7.23)$$

$$G_{\min} \leq G[n] \leq G_{\max} \quad (7.24)$$

with the ratios  $\Delta_{\max} / \Delta_{\min}$  and  $G_{\max} / G_{\min}$  determining the dynamic range of the system. If our objective is to obtain a relatively constant SNR over a range of 40 dB, these ratios can be 100.

*Feedforward adaptation* schemes require us to transmit, in addition to the quantized signal, either the step size  $\Delta[n]$  or the gain  $G[n]$ . Because these values evolve slowly with time, they can be sampled and quantized at a low rate. The overall rate will be the sum of the



bit rate required to transmit the quantized signal plus the bit rate required to transmit either the gain or the step size.

Another class of adaptive quantizers use *feedback adaptation* to avoid having to send information about the step size or gain. In this case, the step size and gain are estimated from the quantizer output, so that they can be recreated at the decoder without any extra information. The corresponding short-time energy can then be estimated through a first-order IIR filter as in Eq. (7.21) or a rectangular window as in Eq. (7.22), but replacing  $x^2[n]$  by  $\hat{x}^2[n]$ .

Another option is to adapt the step size

$$\Delta[n] = P\Delta[n-1] \quad (7.25)$$

where  $P > 1$  if the previous codeword corresponds to the largest positive or negative quantizer level, and  $P < 1$  if the previous codeword corresponds to the smallest positive or negative quantizer level. A similar process can be done for the gain.

APCM exhibits an improvement between 4–8 dB over  $\mu$ -law PCM for the same bit rate.

#### 7.2.4. Differential Quantization

Speech coding is about finding redundancy in the signal and removing it. We know that there is considerable correlation between adjacent samples, because on the average the signal doesn't change rapidly from sample to sample. A simple way of capturing this is to quantize the difference  $d[n]$  between the current sample  $x[n]$  and its predicted value  $\tilde{x}[n]$

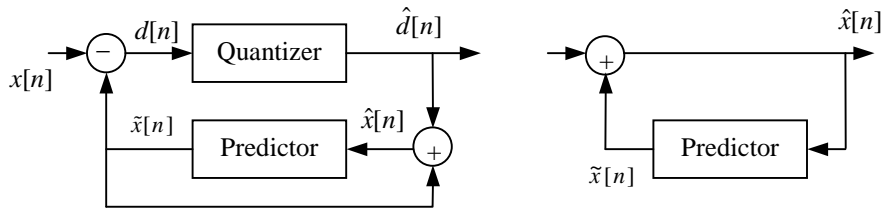
$$d[n] = x[n] - \tilde{x}[n] \quad (7.26)$$

with its quantized value represented as

$$\hat{d}[n] = Q\{d[n]\} = d[n] + e[n] \quad (7.27)$$

where  $e[n]$  is the quantization error. Then, the quantized signal is the sum of the predicted signal  $\tilde{x}[n]$  and the quantized difference  $\hat{d}[n]$

$$\hat{x}[n] = \tilde{x}[n] + \hat{d}[n] = x[n] + e[n] \quad (7.28)$$



**Figure 7.3** Block diagram of a DPCM encoder and decoder with feedback prediction.

If the prediction is good, Eq. (7.28) tells us that the quantization error will be small. Statistically, we need the variance of  $e[n]$  to be lower than that of  $x[n]$  for differential coding to provide any gain. Systems of this type are generically called *Differential Pulse Code Modulation* (DPCM) [11] and can be seen in Figure 7.3.

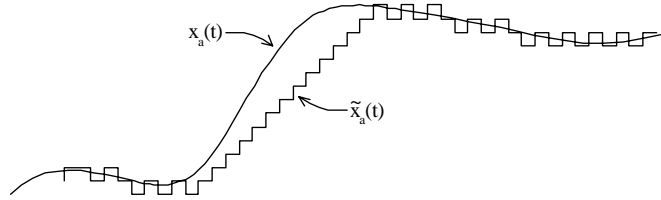
*Delta Modulation* (DM) [47] is a 1-bit DPCM, which predicts the current sample to be the same as the past sample:

$$\tilde{x}[n] = x[n-1] \quad (7.29)$$

so that we transmit whether the current sample is above or below the previous sample.

$$d[n] = \begin{cases} \Delta & x[n] > x[n-1] \\ -\Delta & x[n] \leq x[n-1] \end{cases} \quad (7.30)$$

with  $\Delta$  being the step size. If  $\Delta$  is too small, the reconstructed signal will not increase as fast as the original signal, a condition known as *slope overload distortion*. When the slope is small, the step size  $\Delta$  also determines the peak error; this is known as *granular noise*. Both quantization errors can be seen in Figure 7.4. The choice of  $\Delta$  that minimizes the mean squared error will be a tradeoff between slope overload and granular noise.

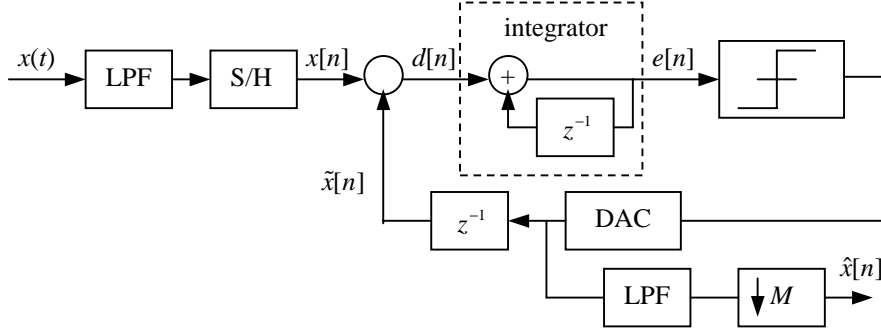


**Figure 7.4** An example of slope overload distortion and granular noise in a DM encoder.

If the signal is oversampled by a factor  $N$ , and the step size is reduced by the same amount (i.e.,  $\Delta/N$ ), the slope overload will be the same, but the granular noise will decrease by a factor  $N$ . While the coder is indeed very simple, sampling rates of over 200 kbps are needed for SNRs comparable to PCM, so DM is rarely used as a speech coder.

However, delta modulation is useful in the design of analog-digital converters, in a variant called sigma-delta modulation [44] shown in Figure 7.5. First the signal is lowpass filtered with a simple analog filter, and then it is oversampled. Whenever the predicted signal  $\tilde{x}[n]$  is below the original signal  $x[n]$ , the difference  $d[n]$  is positive. This difference  $d[n]$  is averaged over time with a digital integrator whose output is  $e[n]$ . If this situation persists, the accumulated error  $e[n]$  will exceed a positive value  $A$ , which causes a 1 to be encoded into the stream  $q[n]$ . A digital-analog converter is used in the loop which increments by one the value of the predicted signal  $\tilde{x}[n]$ . The system acts in the opposite way if the predicted signal  $\tilde{x}[n]$  is above the original signal  $x[n]$  for an extended period of time. Since the signal is oversampled, it changes very slowly from one sample to the next, and this quantization

can be accurate. The advantages of this technique as an analog-digital converter are that inexpensive analog filters can be used and only a simple 1-bit A/D is needed. The signal can next be low-passed filtered with a more accurate digital filter and then downsampled.



**Figure 7.5** A sigma-delta modulator used in an oversampling analog-digital converter.

*Adaptive Delta Modulation* (ADM) combines ideas from adaptive quantization and delta modulation with the so-called *Continuously Variable Slope Delta Modulation* (CVSDM) [22] having a step size that increases

$$\Delta[n] = \begin{cases} \alpha\Delta[n-1] + k_1 & \text{if } e[n], e[n-1] \text{ and } e[n-2] \text{ have same sign} \\ \alpha\Delta[n-1] + k_2 & \text{otherwise} \end{cases} \quad (7.31)$$

with  $0 < \alpha < 1$  and  $0 < k_2 \ll k_1$ . The step size increases if the last three errors have the same sign and decreases otherwise.

Improved DPCM is achieved through linear prediction in which  $\tilde{x}[n]$  is a linear combination of past quantized values  $\hat{x}[n]$

$$\tilde{x}[n] = \sum_{k=1}^p a_k \hat{x}[n-k] \quad (7.32)$$

DPCM systems with fixed prediction coefficients can provide from 4 to 11 dB improvement over direct linear PCM, for prediction orders up to  $p = 4$ , at the expense of increased computational complexity. Larger improvements can be obtained by adapting the prediction coefficients. The coefficients can be transmitted in a feedforward fashion or not transmitted if the feedback scheme is selected.

ADPCM [6] combines differential quantization with adaptive step-size quantization. ITU-T recommendation G.726 uses ADPCM at bit rates of 40, 32, 24, and 16 kbps, with 5, 4, 3, and 2 bits per sample, respectively. It employs an adaptive feedback quantizer and an adaptive feedback pole-zero predictor. Speech at bit rates of 40 and 32 kbps offer toll quality, while the other rates don't. G.727 is called embedded ADPCM because the 2-bit quantizer is embedded into the 3-bit quantizer, which is embedded into the 4-bit quantizer, and into the 5-bit quantizer. This makes it possible for the same codec to use a lower bit rate, with a graceful degradation in quality, if channel capacity is temporarily limited. Earlier

standards G.721 [7, 13] (created in 1984) and G.723 have been subsumed by G.726 and G.727. G.727 has a MOS of 4.1 for 32 kbps and is used in submarine cables. The Windows WAV format also supports a variant of ADPCM. These standards are shown in Table 7.2.

**Table 7.2** Common scalar waveform standards used.

Standard	Bit Rate (kbits/sec)	MOS	Algorithm	Sampling Rate (kHz)
Stereo CD Audio	1411	5.0	16-bit linear PCM	44.1
WAV, AIFF, SND	Variable	-	16/8-bit linear PCM	8, 11.025, 16, 22.05, 44.1, 48
G.711	64	4.3	$\mu$ -law/A-law PCM	8
G.727	40, 32, 24, 16	4.2 (32k)	ADPCM	8
G.722	64, 56, 48		Subband ADPCM	16

Wideband speech (50–7000 Hz) increases intelligibility of fricatives and overall perceived quality. In addition, it provides more subject presence and adds a feeling of transparent communication. ITU-T Recommendation G.722 encodes wideband speech with bit rates of 48, 56, and 64 kbps. Speech is divided into two subbands with QMF filters (see Chapter 5). The upper band is encoded using a 16-kbps ADPCM similar to the G.727 standard. The lower band is encoded using a 48-kbps ADPCM with the 4- and 5-bit quantizers embedded in the 6-bit quantizer. The quality of this system scores almost 1 MOS higher than that of telephone speech.

### 7.3. SCALAR FREQUENCY DOMAIN CODERS

Frequency domain is advantageous because:

1. The samples of a speech signal have a great deal of correlation among them, whereas frequency domain components are approximately uncorrelated and
2. The perceptual effects of masking described in Chapter 2 can be more easily implemented in the frequency domain. These effects are more pronounced for high-bandwidth signals, so frequency-domain coding has been mostly used for CD-quality signals and not for 8-kHz speech signals.

#### 7.3.1. Benefits of Masking

As discussed in Chapter 2, masking is a phenomenon by which human listeners cannot perceive a sound if it is below a certain level. The consequence is that we don't need to encode

such sound. We now illustrate how this masked threshold is computed for MPEG<sup>3</sup>-1 layer 1. Given an input signal  $s[n]$  quantized with  $b$  bits, we obtain the normalized signal  $x[n]$  as

$$x[n] = \frac{s[n]}{N2^{b-1}} \quad (7.33)$$

where  $N = 512$  is the length of the DFT. Then, using a Hanning window,

$$w[n] = 0.5 - 0.5 \cos(2\pi n / N) \quad (7.34)$$

we obtain the log-power spectrum as

$$P[k] = P_0 + 10 \log_{10} \left( \sum_{n=0}^{N-1} w[n] x[n] e^{-j2\pi nk / N} \right) \quad (7.35)$$

where  $P_0$  is the playback SPL, which, in the absence of any volume information, is defined as 90 dB.

*Tonal* components are identified in Eq. (7.35) as local maxima, which exceed neighboring components within a certain bark distance by at least 7 dB. Specifically, bin  $k$  is tonal if and only if

$$P[k] > P[k \pm 1] \quad (7.36)$$

and

$$P[k] > P[k \pm l] + 7 \text{ dB} \quad (7.37)$$

where  $1 < l \leq \Delta_k$ , and  $\Delta_k$  is given by

$$\Delta_k = \begin{cases} 2 & 2 < k < 63 & (170\text{Hz} - 5.5\text{kHz}) \\ 3 & 63 \leq k < 127 & (5.5\text{kHz}, 11\text{kHz}) \\ 6 & 127 \leq k \leq 256 & (11\text{kHz}, 22\text{kHz}) \end{cases} \quad (7.38)$$

so that the power of that tonal masker is computed as the sum of the power in that bin and its left and right adjacent bins:

$$P_{TM}[k] = 10 \log_{10} \left( \sum_{j=-1}^1 10^{0.1P[k+j]} \right) \quad (7.39)$$

The noise maskers are computed from as the sum of power spectrum of the remaining frequency bins  $\bar{k}$  in a critical band not within a neighborhood  $\Delta_k$  of the tonal maskers:

<sup>3</sup> MPEG (Moving Picture Experts Group) is the nickname given to a family of International Standards for coding audiovisual information.

$$P_{NM}[\bar{k}] = 10 \log_{10} \left( \sum_j 10^{0.1P[j]} \right) \quad (7.40)$$

where  $j$  spans a critical band.

To compute the overall masked threshold we need to sum all masking thresholds contributed by each frequency bin  $i$ , which is approximately equal to the maximum (see Chapter 2):

$$T[k] = \max \left( T_h[k], \max_i (T_i[k]) \right) \quad (7.41)$$

In Chapter 2 we saw that whereas temporal postmasking can last from 50 to 300 ms, temporal premasking tends to last about 5 ms. This is also important because when a frequency transform is quantized, the blocking effects of transform's coders can introduce noise above the temporal premasking level that can be audible, since 1024 points corresponds to 23 ms at a 44-kHz sampling rate. To remove this pre-echo distortion, audible in the presence of castanets and other abrupt transient signals, subband filtering has been proposed, whose time constants are well below the 5-ms premasking time constant.

### 7.3.2. Transform Coders

We now use the *Adaptive Spectral Entropy Coding of High Quality Music Signals* (ASPEC) algorithm, which is the basis for the MPEG1 Layer 1 audio coding standard [24], to illustrate how transform coders work. The DFT coefficients are grouped into 128 subbands, and 128 scalar quantizers are used to transmit all the DFT coefficients. It has been empirically found that a difference of less than 1 dB between the original amplitude and the quantized value cannot be perceived. Each subband  $j$  has a quantizer having  $k_j$  levels and step size of  $T_j$  as

$$k_j = 1 + 2 \times \text{rnd} \left( P_j / T_j \right) \quad (7.42)$$

where  $T_j$  is the quantized JND threshold,  $P_j$  is the quantized magnitude of the largest real or imaginary component of the  $j^{\text{th}}$  subband, and  $\text{rnd}(\cdot)$  is the nearest integer rounding function. Entropy coding is used to encode the coefficients of that subband. Both  $T_j$  and  $P_j$  are quantized on a dB scale using 8-bit uniform quantizers with a 170-dB dynamic range, thus with a step size of 0.66 dB. Then they are transmitted as side information.

There are two main methods of obtaining a frequency-domain representation:

1. Through subband filtering via a filterbank (see Chapter 5). When a filterbank is used, the bandwidth of each band is chosen to increase with frequency following a perceptual scale, such as the Bark scale. As shown in Chapter 5, such filterbanks yield perfect reconstruction in the absence of quantization.

2. Through frequency-domain transforms. Instead of using a DFT, higher efficiency can be obtained by the use of an MDCT (see Chapter 5).

The exact implementation of the MPEG1 Layer 1 standard is much more complicated and beyond the scope of this book, though it follows the main ideas described here; the same is true for the popular MPEG1 layer III, also known as MP3. Implementation details can be found in [42].

### 7.3.3. Consumer Audio

Dolby Digital, MPEG, DTS and the Perceptual Audio Coder (PAC) [28] are all audio coders based on frequency-domain coding. Except for MPEG-1, which supports only stereo signals, the rest support multichannel.

Dolby Digital is multichannel digital audio, using lossy AC-3 [54] coding technology from original PCM with a sample rate of 48 kHz at up to 24 bits. The bit rate varies from 64 to 448 kbps, with 384 being the normal rate for 5.1 channels and 192 the normal rate for stereo (with or without surround encoding). Most Dolby Digital decoders support up to 640 kbps. Dolby Digital is the format used for audio tracks on almost all Digital Video/Versatile Discs (DVD). A DVD-5 with only one surround stereo audio stream (at 192 kbps) can hold over 55 hours of audio. A DVD-18 can hold over 200 hours.

MPEG was established in 1988 as part of the joint ISO (International Standardization Organization) / IEC (International Electrotechnical Commission) Technical Committee on Information technology. MPEG-1 was approved in 1992 and MPEG-2 in 1994. Layers I to III define several specifications that provide better quality at the expense of added complexity. MPEG-1 audio is limited to 384 kbps. MPEG1 Layer III audio [23], also known as MP3, is very popular on the Internet, and many compact players exist.

MPEG-2 Audio, one of the audio formats used in DVD, is multichannel digital audio, using lossy compression from 16-bit linear PCM at 48 kHz. Tests have shown that for nearly all types of speech and music, at a data rate of 192 kbps and over, on a stereo channel, scarcely any difference between original and coded versions was observable (ranking of coded item  $> 4.5$ ), with the original signal needing 1.4 Mbps on a CD (reduction by a factor of 7). One advantage of the MPEG Audio technique is that future findings regarding psychoacoustic effects can be incorporated later, so it can be expected that today's quality level using 192 kbps will be achievable at lower data rates in the future. A variable bit rate of 32 to 912 kbps is supported for DVDs.

DTS (Digital Theater Systems) Digital Surround is another multi-channel (5.1) digital audio format, using lossy compression derived from 20-bit linear PCM at 48 kHz. The compressed data rate varies from 64 to 1536 kbps, with typical rates of 768 and 1536 kbps.

### 7.3.4. Digital Audio Broadcasting (DAB)

*Digital Audio Broadcasting* (DAB) is a means of providing current AM and FM listeners with a new service that offers: sound quality comparable to that of compact discs, increased

service availability (especially for reception in moving vehicles), flexible coverage scenarios, and high spectrum efficiency.

Different approaches have been considered for providing listeners with such a service. Currently, the most advanced system is one commonly referred to as Eureka 147 DAB, which has been under development in Europe under the Eureka Project EU147 since 1988. Other approaches include various American in-band systems (IBOC, IBAC, IBRC, FMDigital, and FMeX) still in development, as well as various other systems promising satellite delivery, such as WorldSpace and CD Radio, still in development as well. One satellite-delivery system called MediaStar (formerly Archimedes) proposes to use the Eureka 147 DAB signal structure, such that a single receiver could access both terrestrial and satellite broadcasts.

DAB has been under development since 1981 at the Institut für Rundfunktechnik (IRT) and since 1987 as part of a European research project (Eureka 147). The Eureka 147 DAB specification was standardized by the European Telecommunications Standards Institute (ETSI) in February 1995 as document ETS 300 401, with a draft second edition issued in June 1996. In December 1994, the International Telecommunication Union—Radiocommunication (ITU-R) recommended that this technology, referred to as Digital System A, be used for implementing DAB services.

The Eureka 147 DAB signal consists of multiple carriers within a 1.536-MHz channel bandwidth. Four possible modes of operation define the channel coding configuration, specifying the total number of carriers, the carrier spacing, and also the guard interval duration. Each channel provides a raw data rate of 2304 kbps; after error protection, a useful data rate of anywhere between approximately 600 kbps up to 1800 kbps is available to the service provider, depending on the user-specified multiplex configuration. This useful data rate can be divided into an infinite number of possible configurations of audio and data programs. All audio programs are individually compressed using MUSICAM (MPEG-1 Layer II).

For each useful bit,  $1\frac{1}{3}$  ... 4 bits are transmitted. This extensive redundancy makes it possible to reconstruct the transmitted bit sequence in the receiver, even if part of it is disrupted during transmission (FEC—forward error correction). In the receiver, error concealment can be carried out at the audio reproduction stage, so that residual transmission errors which could not be corrected do not always cause disruptive noise.

## 7.4. CODE EXCITED LINEAR PREDICTION (CELP)

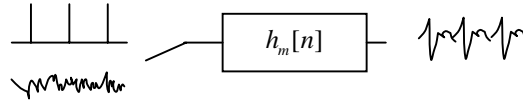
The use of linear predictors removes redundancy in the signal, so that coding of the residual signal can be done with simpler quantizers. We first introduce the LPC vocoder and then introduce coding of the residual signal with a very popular technique called CELP.

### 7.4.1. LPC Vocoder

A typical model for speech production is shown in Figure 7.6, which has a source, or excitation, driving a linear time-varying filter. For voiced speech, the excitation is an impulse train spaced  $P$  samples apart. For unvoiced speech, the source is white random noise. The filter



$h_m[n]$  for frame  $m$  changes at regular intervals, say every 10 ms. If this filter is represented with linear predictive coding, it is called an *LPC vocoder* [3].



**Figure 7.6** Block diagram of an LPC vocoder.

In addition to transmitting the gain and LPC coefficients, the encoder has to determine whether the frame is voiced or unvoiced, as well as the pitch period  $P$  for voiced frames.

The LPC vocoder produces reasonable quality for unvoiced frames, but often results in somewhat mechanical sound for voiced sounds, and a buzzy quality for voiced fricatives. More importantly, the LPC vocoder is quite sensitive to voicing and pitch errors, so that an accurate pitch tracker is needed for reasonable quality. The LPC vocoder also performs poorly in the presence of background noise. Nonetheless, it can be highly intelligible. The Federal Standard 1015 [55], proposed for secure communications, is based on a 2.4-kbps LPC vocoder.

It's also possible to use linear predictive coding techniques together with Huffman coding [45] to achieve lossless compression of up to 50%.

### 7.4.2. Analysis by Synthesis

*Code Excited Linear Prediction* (CELP) [5] is an umbrella for a family of techniques that quantize the LPC residual using VQ, thus the term *code excited*, using analysis by synthesis. In addition CELP uses the fact that the residual of voiced speech has periodicity and can be used to predict the residual of the current frame. In CELP coding the LPC coefficients are quantized and transmitted (feedforward prediction), as well as the codeword index. The prediction using LPC coefficients is called *short-term prediction*. The prediction of the residual based on pitch is called *long-term prediction*. To compute the quantized coefficients we use an *analysis-by-synthesis* technique, which consists of choosing the combination of parameters whose reconstructed signal is closest to the analysis signal. In practice, not all coefficients of a CELP coder are estimated in an analysis-by-synthesis manner.

We first estimate the  $p^{\text{th}}$ -order LPC coefficients from the samples  $x[n]$  for frame  $t$  using the autocorrelation method, for example. We then quantize the LPC coefficients to  $(a_1, a_2, \dots, a_p)$  with the techniques described in Section 7.4.5. The residual signal  $e[n]$  is obtained by inverse filtering  $x[n]$  with the quantized LPC filter

$$e[n] = \sum_{i=1}^p a_i x[n-i] \quad (7.43)$$

Given the transfer function of the LPC filter

$$H(z) = \frac{1}{A(z)} = \frac{1}{1 - \sum_{i=1}^p a_i z^{-i}} = \sum_{i=0}^{\infty} h_i z^{-i} \quad (7.44)$$

we can obtain the first  $M$  coefficients of the impulse response  $h[n]$  of the LPC filter by driving it with an impulse as

$$h[n] = \begin{cases} 1 & n = 0 \\ \sum_{i=1}^n a_i h[n-i] & 0 < n < p \\ \sum_{i=1}^p a_i h[n-i] & p \leq n < M \end{cases} \quad (7.45)$$

so that if we quantize a frame of  $M$  samples of the residual  $\mathbf{e} = (e[0], e[1], \dots, e[M-1])^T$  to  $\mathbf{e}_i = (e_i[0], e_i[1], \dots, e_i[M-1])^T$ , we can compute the reconstructed signal  $\hat{x}_i[n]$  as

$$\hat{x}_i[n] = \sum_{m=0}^n h[m] e_i[n-m] + \sum_{m=n+1}^{\infty} h[m] e[n-m] \quad (7.46)$$

where the second term in the sum depends on the residual for previous frames, which we already have. Let's define signal  $r_0[n]$  as the second term of Eq. (7.46):

$$r_0[n] = \sum_{m=n+1}^{\infty} h[m] e[n-m] \quad (7.47)$$

which is the output of the LPC filter when there is no excitation for frame  $t$ . The important thing to note is that  $r_0[n]$  does not depend on  $e_i[n]$

It is convenient to express Eqs. (7.46) and (7.47) in matrix form as

$$\hat{\mathbf{x}}_i = \mathbf{H} \mathbf{e}_i + \mathbf{r}_0 \quad (7.48)$$

where matrix  $\mathbf{H}$  corresponds to the LPC filtering operation with its memory set to 0:

$$\mathbf{H} = \begin{bmatrix} h_0 & 0 & \cdots & 0 & 0 \\ h_1 & h_0 & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ h_{M-1} & h_{M-2} & \cdots & h_0 & 0 \\ h_M & h_{M-1} & \cdots & h_1 & h_0 \end{bmatrix} \quad (7.49)$$

Given the large dynamic range of the residual signal, we use gain-shape quantization, where we quantize the gain and the gain-normalized residual separately:

$$\mathbf{e}_i = \lambda \mathbf{c}_i \quad (7.50)$$

where  $\lambda$  is the gain and  $\mathbf{c}_i$  is the codebook entry  $i$ . This codebook is known as the *fixed codebook* because its vectors do not change from frame to frame. Usually the size of the codebook is selected as  $2^N$  so that full use is made of all  $N$  bits. Codebook sizes typically vary from 128 to 1024. Combining Eq. (7.48) with Eq. (7.50), we obtain

$$\hat{\mathbf{x}}_i = \lambda \mathbf{H} \mathbf{c}_i + \mathbf{r}_0 \quad (7.51)$$

The error between the original signal  $\mathbf{x}$  and the reconstructed signal  $\hat{\mathbf{x}}_i$  is

$$\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}_i \quad (7.52)$$

The optimal gain  $\lambda$  and codeword index  $i$  are the ones that minimize the squared error between the original signal and the reconstructed<sup>4</sup> signal:

$$E(i, \lambda) = \|\mathbf{x} - \hat{\mathbf{x}}_i\|^2 = \|\mathbf{x} - \lambda \mathbf{H} \mathbf{c}_i - \mathbf{r}_0\|^2 = \|\mathbf{x} - \mathbf{r}_0\|^2 + \lambda^2 \mathbf{c}_i^T \mathbf{H}^T \mathbf{H} \mathbf{c}_i - 2\lambda \mathbf{c}_i^T \mathbf{H}^T (\mathbf{x} - \mathbf{r}_0) \quad (7.53)$$

where the term  $\|\mathbf{x} - \mathbf{r}_0\|^2$  does not depend on  $\lambda$  or  $i$  and can be neglected in the minimization. For a given  $\mathbf{c}_i$ , the gain  $\lambda_i$  that minimizes Eq. (7.53) is given by

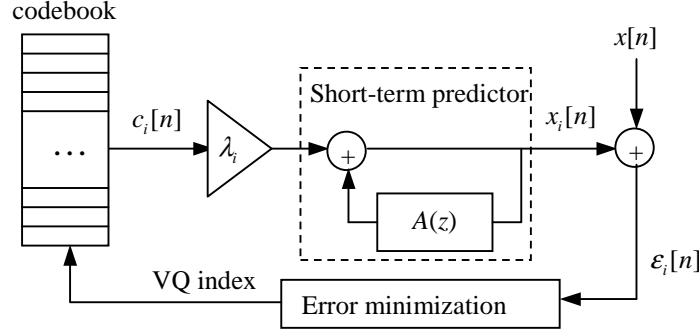
$$\lambda_i = \frac{\mathbf{c}_i^T \mathbf{H}^T (\mathbf{x} - \mathbf{r}_0)}{\mathbf{c}_i^T \mathbf{H}^T \mathbf{H} \mathbf{c}_i} \quad (7.54)$$

Inserting Eq. (7.54) into (7.53) lets us compute the index  $j$  as the one that minimizes

$$j = \arg \min_i \left\{ -\frac{(\mathbf{c}_i^T \mathbf{H}^T (\mathbf{x} - \mathbf{r}_0))^2}{\mathbf{c}_i^T \mathbf{H}^T \mathbf{H} \mathbf{c}_i} \right\} \quad (7.55)$$

So we first obtain the codeword index  $j$  according to Eq. (7.55) and then the gain  $\lambda_j$  according to Eq. (7.54), which is scalarly quantized to  $\hat{\lambda}_j$ . Both codeword index  $j$  and  $\hat{\lambda}_j$  are transmitted. In the algorithm described here, we first chose the quantized LPC coefficients  $(a_1, a_2, \dots, a_p)$  independently of the gains and codeword index, and then we chose the codeword index independently of the quantized gain  $\hat{\lambda}_j$ . This procedure is called *open-loop* estimation, because some parameters are obtained independently of the others. This is shown in Figure 7.7. *Closed-loop* estimation [49] means that all possible combinations of quantized parameters are explored. Closed-loop is more computationally expensive but yields lower squared error.

<sup>4</sup> A beginner's mistake is to find the codebook index that minimizes the squared error of the residual. This does not minimize the difference between the original signal and the reconstructed signal.



**Figure 7.7** Analysis-by-synthesis principle used in a basic CELP.

### 7.4.3. Pitch Prediction: Adaptive Codebook

The fact that speech is highly periodic during voiced segments can also be used to reduce redundancy in the signal. This can be done by predicting the residual signal  $e[n]$  at the current vector with samples from the past residual signal shifted a pitch period  $t$ :

$$e[n] = \lambda_i^a e[n-t] + \lambda_i^f c_i^f[n] = \lambda_i^a c_i^a[n] + \lambda_i^f c_i^f[n] \quad (7.56)$$

Using the matrix framework we described before, Eq. (7.56) can be expressed as

$$\mathbf{e}_{ii} = \lambda_i^a \mathbf{c}_i^a + \lambda_i^f \mathbf{c}_i^f \quad (7.57)$$

where we have made use of an *adaptive codebook* [31], where  $\mathbf{c}_i^a$  is the adaptive codebook entry  $j$  with corresponding gain  $\lambda^a$ , and  $\mathbf{c}_i^f$  is the fixed or stochastic codebook entry  $i$  with corresponding gain  $\lambda^f$ . The adaptive codebook entries are segments of the recently synthesized excitation signal

$$\mathbf{c}_i^a = (e[-t], e[1-t], \dots, e[M-1-t])^T \quad (7.58)$$

where  $t$  is the delay which specifies the start of the adaptive codebook entry  $t$ . The range of  $t$  is often between 20 and 147, since this can be encoded with 7 bits. This corresponds to a range in pitch frequency between 54 and 400 Hz for a sampling rate of 8 kHz.

The contribution of the adaptive codebook is much larger than that of the stochastic codebook for voiced sounds. So we generally search for the adaptive codebook first, using Eq. (7.58) and a modified version of Eqs. (7.55), (7.54), replacing  $i$  by  $t$ . Closed-loop search of both  $t$  and gain here often yields a much larger error reduction.

#### 7.4.4. Perceptual Weighting and Postfiltering

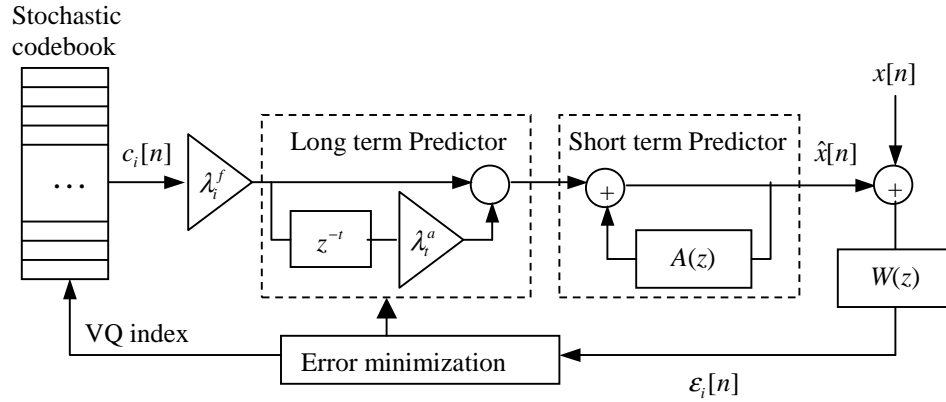
The objective of speech coding is to reduce the bit rate while maintaining a perceived level of quality; thus, minimization of the error is not necessarily the best criterion. A perceptual weighting filter tries to shape the noise so that it gets masked by the speech signal (see Chapter 2). This generally means that most of the quantization noise energy is located in spectral regions where the speech signal has most of its energy. A common technique [4] consists in approximating this perceptual weighting with a linear filter

$$W(z) = \frac{A(z/\beta)}{A(z/\gamma)} \quad (7.59)$$

where  $A(z)$  is the predictor polynomial

$$A(z) = 1 - \sum_{i=1}^p a_i z^{-i} \quad (7.60)$$

Choosing  $\gamma$  and  $\beta$  so that  $0 < \gamma < \beta \leq 1$ , implies that the roots of  $A(z/\beta)$  and  $A(z/\gamma)$  will move closer to the origin of the unit circle than the roots of  $A(z)$ , thus resulting in a frequency response with wider resonances. This perceptual filter therefore deemphasizes the contribution of the quantization error near the formants. A common choice of parameters is  $\beta = 1.0$  and  $\gamma = 0.8$ , since it simplifies the implementation. This filter can easily be included in the matrix  $\mathbf{H}$ , and a CELP coder incorporating the perceptual weighting is shown in Figure 7.8.



**Figure 7.8** Diagram of a CELP coder. Both long-term and short-term predictors are used, together with a perceptual weighting.

Despite the perceptual weighting filter, the reconstructed signal still contains audible noise. This filter reduces the noise in those frequency regions that are perceptually irrelevant without degrading the speech signal. The postfilter generally consists of a short-term postfilter to emphasize the formant structure and a long-term postfilter to enhance the periodicity

of the signal [10]. One possible implementation follows Eq. (7.59) with values of  $\beta = 0.5$  and  $\gamma = 0.75$ .

#### 7.4.5. Parameter Quantization

To achieve a low bit rate, all the coefficients need to be quantized. Because of its coding efficiency, vector quantization is the compression technique of choice to quantize the predictor coefficients. The LPC coefficients cannot be quantized directly, because small errors produced in the quantization process may result in large changes in the spectrum and possibly unstable filters. Thus, equivalent representations that guarantee stability are used, such as reflection coefficients, log-area ratios, and the line spectral frequencies (LSF) described in Chapter 6. LSF are used most often, because it has been found empirically that they behave well when they are quantized and interpolated [2]. For 8 kHz, 10 predictor coefficients are often used, which makes using a single codebook impractical because of the large dimension of the vector. Split-VQ [43] is a common choice, where the vectors are divided into several subvectors, and each is vector quantized. Matrix quantization can also be used to exploit the correlation of these subvectors across consecutive time frames. *Transparent quality*, defined as average spectral distortion below 1 dB with no frames above 4 dB, can be achieved with fewer than 25 bits per frame.

A frame typically contains around 20 to 30 milliseconds, which at 8 kHz represents 160–240 samples. Because of the large vector dimension, it is impractical to quantize a whole frame with a single codebook. To reduce the dimensionality, the frame is divided into four or more nonoverlapping sub-frames. The LSF coefficients for each subframe are linearly interpolated between the two neighboring frames.

A typical range of the pitch prediction for an 8-kHz sampling rate goes from 2 to 20 ms, from 20 to 147 samples, 2.5 ms to 18.375 ms, which can be encoded with 7 bits. An additional bit is often used to encode fractional delays for the lower pitch periods. These fractional delays can be implemented through upsampling as described in Chapter 5. The subframe gain of the adaptive codebook can be effectively encoded with 3 or 4 bits. Alternatively, the gains of all sub-frames within a frame can be encoded through VQ, resulting in more efficient compression.

The fixed codebook can be trained from data using the techniques described in Chapter 4. This will offer the lowest distortion for the training set but doesn't guarantee low distortion for mismatched test signals. Also, it requires additional storage, and full search increases computation substantially.

Since subframes should be approximately white, the codebook can be populated from samples of a white process. A way of reducing computation is to let those noise samples be only +1, 0, or -1, because only additions are required. Codebooks of a specific type, known as *algebraic codebooks* [1], offer even more computational savings because they contain many 0s. Locations for the 4 pulses per subframe under the G.729 standard are shown in Table 7.3.

Full search can efficiently be done with this codebook structure. Algebraic codebooks can provide almost as low distortion as trained codebooks can, with low computational complexity.

**Table 7.3** Algebraic codebooks for the G.729 standard. Each of the four codebooks has one pulse in one possible location indicated by 3 bits for the first three codebooks and 4 bits for the last codebook. The sign is indicated by an additional bit. A total of 17 bits are needed to encode a 40-sample subframe.

Amplitude	Positions
$\pm 1$	0, 5, 10, 15, 20, 25, 30, 35
$\pm 1$	1, 6, 11, 16, 21, 26, 31, 36
$\pm 1$	2, 7, 12, 17, 22, 27, 32, 37
$\pm 1$	3, 8, 13, 18, 23, 28, 33, 38
	4, 9, 14, 19, 24, 29, 34, 39

#### 7.4.6. CELP Standards

There are many standards for speech coding based on CELP, offering various points in the bit-rate/quality plane, mostly depending on when they were created and how refined the technology was at that time.

Voice over Internet Protocol (Voice over IP) consists of transmission of voice through data networks such as the Internet. H.323 is an umbrella standard which references many other ITU-T recommendations. H.323 provides the system and component descriptions, call model descriptions, and call signaling procedures. For audio coding, G.711 is mandatory, while G.722, G.728, G.723.1, and G.729 are optional. G.728 is a low-delay CELP coder that offers toll quality at 16 kbps [9], using a feedback 50<sup>th</sup>-order predictor, but no pitch prediction. G.729 [46] offers toll quality at 8 kbps, with a delay of 10 ms. G.723.1, developed by DSP Group, including Audiocodes Ltd., France Telecom, and the University of Sherbrooke, has slightly lower quality at 5.3 and 6.3 kbps, but with a delay of 30 ms. These standards are shown in Table 7.4.

**Table 7.4** Several CELP standards used in the H.323 specification used for teleconferencing and voice streaming through the internet.

Standard	Bit Rate (kbps)	MOS	Algorithm	H.323	Comments
G.728	16	4.0	No pitch prediction	Optional	Low -delay
G.729	8	3.9	ACELP	Optional	
G.723.1	5.3, 6.3	3.9	ACELP for 5.3k	Optional	

In 1982, the Conference of European Posts and Telegraphs (CEPT) formed a study group called the Groupe Spécial Mobile (GSM) to study and develop a pan-European public land mobile system. In 1989, GSM responsibility was transferred to the European Telecommunication Standards Institute (ETSI), and the phase I GSM specifications were published in 1990. Commercial service was started in mid 1991, and by 1993 there were 36 GSM networks in 22 countries, with 25 additional countries considering or having already selected GSM. This is not only a European standard; South Africa, Australia, and many Middle and Far East countries have chosen GSM. The acronym GSM now stands for Global System for Mobile telecommunications. The GSM group studied several voice coding algorithms on the basis of subjective speech quality and complexity (which is related to cost, processing delay, and power consumption once implemented) before arriving at the choice of a Regular Pulse Excited–Linear Predictive Coder (RPE-LPC) with a Long Term Predictor loop [56]. Neither the original full-rate at 13 kbps [56] nor the half-rate at 5.6 kbps [19] achieves toll quality, though the enhanced full-rate (EFR) standard based on ACELP [26] has toll quality at the same rates.

The *Telecommunication Industry Association* (TIA) and the *Electronic Industries Alliance* (EIA) are organizations accredited by the *American National Standards Institute* (ANSI) to develop voluntary industry standards for a wide variety of telecommunication products. TR-45 is the working group within TIA devoted to mobile and personal communication systems. Time Division Multiple Access (TDMA) is a digital wireless technology that divides a narrow radio channel into framed time slots (typically 3 or 8) and allocates a slot to each user. The TDMA Interim Standard 54, or TIA/EIA/IS54, was released in early 1991 by both TIA and EIA. It is available in North America at both the 800-MHz and 1900-MHz bands. IS54 [18] at 7.95 kbps is used in North America's TDMA (Time Division Multiple Access) digital telephony and has quality similar to the original full-rate GSM. TDMA IS-136 is an update released in 1994.

**Table 7.5** CELP standards used in cellular telephony.

Standard	Bit Rate (kbps)	MOS	Algorithm	Cellular	Comments
Full-rate GSM	13	3.6	VSELP RTE-LTP	GSM	
EFR GSM	12.2	4.5	ACELP	GSM	
IS-641	7.4	4.1	ACELP	PCS1900	
IS-54	7.95	3.9	VSELP	TDMA	
IS-96a	max 8.5	3.9	QCELP	CDMA	Variable-rate

*Code Division Multiple Access* (CDMA) is a form of *spread spectrum*, a family of digital communication techniques that have been used in military applications for many years. The core principle is the use of noiselike carrier waves, and, as the name implies, bandwidths much wider than that required for simple point-to-point communication at the same data rate. Originally there were two motivations: either to resist enemy efforts to jam



the communications (anti-jam, or AJ) or to hide the fact that communication was even taking place, sometimes called low probability of intercept (LPI). The service started in 1996 in the United States, and by the end of 1999 there were 50 million subscribers worldwide. IS-96 QCELP [14], used in North America's CDMA, offers variable-rate coding at 8.5, 4, 2 and 0.8 kbps. The lower bit rate is transmitted when the coder detects background noise. TIA/EIA/IS-127-2 is a standard for an enhanced variable-rate codec, whereas TIA/EIA/IS-733-1 is a standard for high-rate. Standards for CDMA, TDMA, and GSM are shown in Table 7.5.

*Third generation (3G)* is the generic term used for the next generation of mobile communications systems. 3G systems will provide enhanced services to those—such as voice, text, and data—predominantly available today. The Universal Mobile Telecommunications System (UMTS) is a part of ITU's International Mobile Telecommunications (IMT)-2000 vision of a global family of third-generation mobile communications systems. It has been assigned to the frequency bands 1885–2025 and 2110–2200 MHz. The first networks are planned to launch in Japan in 2001, with European countries following in early 2002. A major part of 3G is General Packet Radio Service (GPRS), under which carriers charge by the packet rather than by the minute. The speech coding standard for CDMA2000, the umbrella name for the third-generation standard in the United States, is expected to gain approval late in 2000. An adaptive multi rate wideband speech codec has also been proposed for the GSM's 3G [16], which has five modes of operation from 24 kbps down to 9.1 kbps.

While most of the work described above uses a sampling rate of 8 kHz, there has been growing interest in using CELP techniques for high bandwidth and particularly in a scalable way so that a basic layer contains the lower frequency and the higher layer either is a full-band codec [33] or uses a parametric model [37].

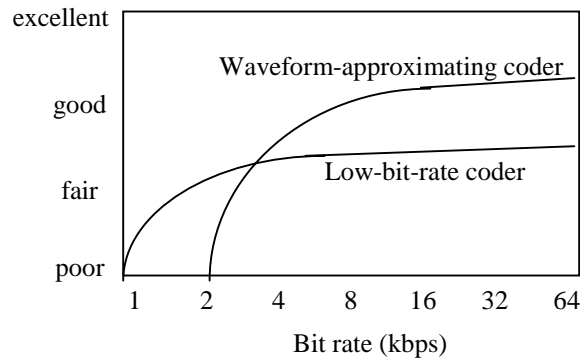
## 7.5. LOW-BIT RATE SPEECH CODERS

In this section we describe a number of low-bit-rate speech coding techniques including the mixed-excitation LPC vocoder, harmonic coding, and waveform interpolation. These coding techniques are also used extensively in speech synthesis.

*Waveform-approximating coders* are designed to minimize the difference between the original signal and the coded signal. Therefore, they produce a reconstructed signal whose SNR goes to infinity as the bit rate increases, and they also behave well when the input signal is noisy or music. In this category we have the scalar waveform coders of Section 7.2, the frequency-domain coders of Section 7.3, and the CELP coders of Section 7.4.

Low-bit-rate coders, on the other hand, do not attempt to minimize the difference between the original signal and the quantized signal. Since these coders are designed to operate at low bit rates, their SNR does not generally approach infinity even if a large bit rate is used. The objective is to compress the original signal with another one that is perceptually equivalent. Because of the reliance on an inaccurate model, these low-bit-rate coders often distort the speech signal even if the parameters are not quantized. In this case, the distortion can consist of more than quantization noise. Furthermore, these coders are more sensitive to the presence of noise in the signal, and they do not perform as well on music.

In Figure 7.9 we compare the MOS of waveform approximating coders and low-bit-rate coders as a function of the bit rate. CELP uses a model of speech to obtain as much prediction as possible, yet allows for the model not to be exact, and thus is a waveform-approximating coder. CELP is a robust coder that works reasonably well when the assumption of only a clean speech signal breaks either because of additive noise or because there is music in the background. Researchers are working on the challenging problem of creating more scalable coders that offer best performance at all bit rates.



**Figure 7.9** Typical subjective performance of waveform-approximating and low-bit-rate coders as a function of the bit rate. Note that waveform-approximating coders are a better choice for bit rates higher than about 3 kbps, whereas parametric coders are a better choice for lower bit rates. The exact cutoff point depends on the specific algorithms compared.

### 7.5.1. Mixed-Excitation LPC Vocoder

The main weakness of the LPC vocoder is the binary decision between voiced and unvoiced speech, which results in errors especially for noisy speech and voiced fricatives. By having a separate voicing decision for each of a number of frequency bands, the performance can be enhanced significantly [38]. The new proposed U.S. Federal Standard at 2.4 kbps is a Mixed Excitation Linear Prediction (MELP) LPC vocoder [39], which has a MOS of about 3.3. This exceeds the quality of the older 4800-bps federal standard 1016 [8] based on CELP. The bit rate of the proposed standard can be reduced while maintaining the same quality by jointly quantizing several frames together [57]. A hybrid codec that uses MELP in strongly voiced regions and CELP in weakly voiced and unvoiced regions [53] has shown to yield lower bit rates. MELP can also be combined with the waveform interpolation technique of Section 7.5.3 [50].

### 7.5.2. Harmonic Coding

Sinusoidal coding decomposes the speech signal [35] or the LP residual signal [48] into a sum of sinusoids. The case where these sinusoids are harmonically related is of special in-

terest for speech synthesis (see Chapter 16), so we will concentrate on it in this section, even though a similar treatment can be followed for the case where the sinusoids are not harmonically related. In fact, a combination of harmonically related and nonharmonically related sinusoids can also be used [17]. We show in Section 7.5.2.2 that we don't need to transmit the phase of the sinusoids, only the magnitude.

As shown in Chapter 5, a periodic signal  $\tilde{s}[n]$  with period  $T_0$  can be expressed as a sum of  $T_0$  harmonic sinusoids

$$\tilde{s}[n] = \sum_{l=0}^{T_0-1} A_l \cos(nl\omega_0 + \phi_l) \quad (7.61)$$

whose frequencies are multiples of the fundamental frequency  $\omega_0 = 2\pi/T_0$ , and where  $A_l$  and  $\phi_l$  are the sinusoid amplitudes and phases, respectively. If the pitch period  $T_0$  has fractional samples, the sum in Eq. (7.61) includes only the integer part of  $T_0$  in the summation. Since a real signal  $s[n]$  will not be perfectly periodic in general, we have a modeling error

$$e[n] = s[n] - \tilde{s}[n] \quad (7.62)$$

We can use short-term analysis to estimate these parameters from the input signal  $s[n]$  at frame  $k$ , in the neighborhood of  $t = kN$ , where  $N$  is the frame shift:

$$s_k[n] = s[n]w_k[n] = s[n]w[kN - n] \quad (7.63)$$

if we make the assumption that the sinusoid parameters for frame  $k$  ( $\omega_0^k$ ,  $A_l^k$  and  $\phi_l^k$ ) are constant within the frame.

At resynthesis time, there will be discontinuities at unit boundaries, due to the block processing, unless we specifically smooth the parameters over time. One way of doing this is with overlap-add method between frames  $(k-1)$  and  $k$ :

$$\hat{s}[n] = w[n]\tilde{s}^{k-1}[n] + w[n-N]\tilde{s}^k[n-N] \quad (7.64)$$

where the window  $w[n]$  must be such that

$$w[n] + w[n-N] = 1 \quad (7.65)$$

to achieve perfect reconstruction. This is the case for the common Hamming and Hanning windows.

This harmonic model [35] is similar to the classic filterbank, though rather than the whole spectrum we transmit only the fundamental frequency  $\omega_0$  and the amplitudes  $A_l$  and phases  $\phi_l$  of the harmonics. This reduced representation doesn't result in loss of quality for a frame shift  $N$  that corresponds to 12 ms or less. For unvoiced speech, using a default pitch of 100 Hz results in acceptable quality.

### 7.5.2.1. Parameter Estimation

For simplicity in the calculations, let's define  $\tilde{s}[n]$  as a sum of complex exponentials

$$\tilde{s}[n] = \sum_{l=0}^{T_0-1} A_l \exp\{j(nl\omega_0 + \phi_l)\} \quad (7.66)$$

and perform short-time Fourier transform with a window  $w[n]$

$$\tilde{S}_w(\omega) = \sum_{l=0}^{T_0-1} A_l e^{j\phi_l} W(\omega - l\omega_0) \quad (7.67)$$

where  $W(\omega)$  is the Fourier transform of the window function. The goal is to estimate the sinusoid parameters as those that minimize the squared error:

$$E = |S(\omega) - \tilde{S}_w(\omega)|^2 \quad (7.68)$$

If the main lobes of the analysis window do not overlap, we can estimate the phases  $\phi_l$  as

$$\phi_l = \arg S(l\omega_0) \quad (7.69)$$

and the amplitudes  $A_l$  as

$$A_l = \frac{|S(l\omega_0)|}{W(0)} \quad (7.70)$$

For example, the Fourier transform of a  $(2N + 1)$  point rectangular window centered around the origin is given by

$$W(\omega) = \frac{\sin((2N + 1)\omega/2)}{\sin(\omega/2)} \quad (7.71)$$

whose main lobes will not overlap in Eq. (7.67) if  $2T_0 < 2N + 1$ : i.e., the window contains at least two pitch periods. The implicit assumption in the estimates of Eqs. (7.69) and (7.70) is that there is no spectral leakage, but a rectangular window does have significant spectral leakage, so a different window is often used in practice. For windows such as Hanning or Hamming, which reduce the leakage significantly, it has been found experimentally that these estimates are correct if the window contains at least two and a half pitch periods.

Typically, the window is centered around 0 (nonzero in the interval  $-N \leq n \leq N$ ) to avoid numerical errors in estimating the phases.

Another implicit assumption in Eqs. (7.69) and (7.70) is that we know the fundamental frequency  $\omega_0$  ahead of time. Since, in practice, this is not the case, we can estimate it as the one which minimizes Eq. (7.68). This pitch-estimation method can generate pitch doubling or tripling when a harmonic falls within a formant that accounts for the majority of the signal's energy.

Voiced/unvoiced decisions can be computed from the ratio between the energy of the signal and that of the reconstruction error

$$SNR = \frac{\sum_{n=-N}^N |s[n]|^2}{\sum_{n=-N}^N |s[n] - \tilde{s}[n]|^2} \quad (7.72)$$

where it has been empirically found that frames with SNR higher than 13 dB are generally voiced and lower than 4 dB unvoiced. In between, the signal is considered to contain a mixed excitation. Since speech is not perfectly stationary within the analysis frame, even noise-free periodic signals will yield finite SNR.

For unvoiced speech, a good assumption is to default to a pitch of 100 Hz. The use of fewer sinusoids leads to perceptual artifacts.

Improved quality can be achieved by using an analysis-by-synthesis framework [17, 34] since the closed-loop estimation is more robust to pitch-estimation and voicing decision errors.

### 7.5.2.2. Phase Modeling

An impulse train  $e[n]$ , a periodic excitation, can be expressed as a sum of complex exponentials

$$e[n] = T_0 \sum_{k=-\infty}^{\infty} \delta[n - n_0 - kT_0] = \sum_{l=0}^{T_0-1} e^{j(n-n_0)\omega_0 l} \quad (7.73)$$

which, if passed through a filter  $H(\omega) = A(\omega) \exp \Phi(\omega)$ , will generate

$$s[n] = \sum_{l=0}^{T_0-1} A(l\omega_0) \exp\{j[(n-n_0)\omega_0 l + \Phi(l\omega_0)]\} \quad (7.74)$$

Comparing Eq. (7.66) with (7.74), the phases of our sinusoidal model are given by

$$\phi_l = -n_0 \omega_0 l + \Phi(l\omega_0) \quad (7.75)$$

Since the sinusoidal model has too many parameters to lead to low-rate coding, a common technique is to not encode the phases. In Chapter 6 we show that if a system is considered minimum phase, the phases can be uniquely recovered from knowledge of the magnitude spectrum.

The magnitude spectrum is known at the pitch harmonics, and the remaining values can be filled in by interpolation: e.g., linear or cubic splines [36]. This interpolated magnitude spectrum can be approximated through the real cepstrum:

$$|\tilde{A}(\omega)| = c_0 + 2 \sum_{k=1}^K c_k \cos(k\omega) \quad (7.76)$$

and the phase, assuming a minimum phase system, is given by

$$\tilde{\Phi}(\omega) = -2 \sum_{k=1}^K c_k \sin(k\omega) \quad (7.77)$$

The phase  $\phi_0(t)$  of the first harmonic between frames  $(k-1)$  and  $k$  can be obtained from the instantaneous frequency  $\omega_0(t)$

$$\phi_0(t) = \phi_0((k-1)N) + \int_{(k-1)N}^t \omega_0(t) dt \quad (7.78)$$

if we assume the frequency  $\omega_0(t)$  in that region to vary linearly between frames  $(k-1)$  and  $k$ :

$$\omega_0(t) = \omega_0^{k-1} + \frac{\omega_0^k - \omega_0^{k-1}}{N} t \quad (7.79)$$

and insert Eq. (7.79) into (7.78), evaluating at  $t = kN$ , to obtain

$$\phi_0^k = \phi_0(kN) = \phi_0((k-1)N) + (\omega_0^{k-1} + \omega_0^k)(N/2) \quad (7.80)$$

the phase of the sinusoid at  $\omega_0$  as a function of the fundamental frequencies at frames  $(k-1)$ ,  $k$  and the phase at frame  $(k-1)$ :

$$\phi_l^k = \Phi^k(l\omega_0) + l\phi_0^k \quad (7.81)$$

The phases computed by Eqs. (7.80) and (7.81) are a good approximation in practice for perfectly voiced sounds. For unvoiced sounds, random phases are needed, or else the reconstructed speech sounds buzzy. Voiced fricatives and many voiced sounds have an aspiration component, so that a mixed excitation is needed to represent them. In these cases, the source is split into different frequency bands and each band is classified as either voiced or unvoiced. Sinusoids in voiced bands use the phases described above, whereas sinusoids in unvoiced bands have random phases.

### 7.5.2.3. Parameter Quantization

To quantize the sinusoid amplitudes, we can use an LPC fitting and then quantize the line spectral frequencies. Also we can do a cepstral fit and quantize the cepstral coefficients. To be more effective, a mel scale should be used.

While these approaches help in reducing the number of parameters and in quantizing those parameters, they are not the most effective way of quantizing the sinusoid amplitudes. A technique called *Variable-Dimension Vector Quantization* (VDVQ) [12] has been devised to address this. Each codebook vector  $\mathbf{c}_i$  has a fixed dimension  $N$  determined by the length of the FFT used. The vector of sinusoid amplitudes  $\mathbf{A}$  has a dimension  $l$  that depends on the number of harmonics and thus the pitch of the current frame. To compute the distance between  $\mathbf{A}$  and  $\mathbf{c}_i$ , the codebook vectors are resampled to a size  $l$  and the distance is computed between two vectors of dimension  $l$ . Euclidean distance of the log-amplitudes is often used. In this method, only the distance at the harmonics is evaluated instead of the distance at the points in the envelope that are actually not present in the signal. Also, this technique does

not suffer from inaccuracies of the model used, such as the inability of linear predictive coding to model nasals.

### 7.5.3. Waveform Interpolation

The main idea behind waveform interpolation (WI) [29] is that the pitch pulse changes slowly over time for voiced speech. During voiced segments, the speech signal is nearly periodic. WI coders can operate as low as 2.4 kbps.

Starting at an arbitrary time instant, it is easy to identify a first pitch cycle  $x_1[n]$ , a second  $x_2[n]$ , a third  $x_3[n]$ , and so on. We then express our signal  $x[n]$  as a function of these pitch cycle waveforms  $x_m[n]$

$$x[n] = \sum_{m=-\infty}^{\infty} x_m[n - t_m] \quad (7.82)$$

where  $P_m = t_m - t_{m-1}$  is the pitch period at time  $t_m$  in samples, and the pitch cycle is a windowed version of the input

$$x_m[n] = w_m[n]x[n] \quad (7.83)$$

—for example, with a rectangular window. To transmit the signal in a lossless fashion we need to transmit all pitch waveforms  $x_m[n]$ .

If the signal is perfectly periodic, we need to transmit only one pitch waveform  $x_m[n]$  and the pitch period  $P$ . In practice, voiced signals are not perfectly periodic, so that we need to transmit more than just one pitch waveform. On the other hand, voiced speech is nearly periodic, and consecutive pitch waveforms are very similar. Thus, we probably do not need to transmit all, and we could send every other pitch waveform, for example.

It is convenient to define a two-dimensional surface  $u[n, l]$  (shown in Figure 7.10) such that the pitch waveform  $x_m[n]$  can be obtained as

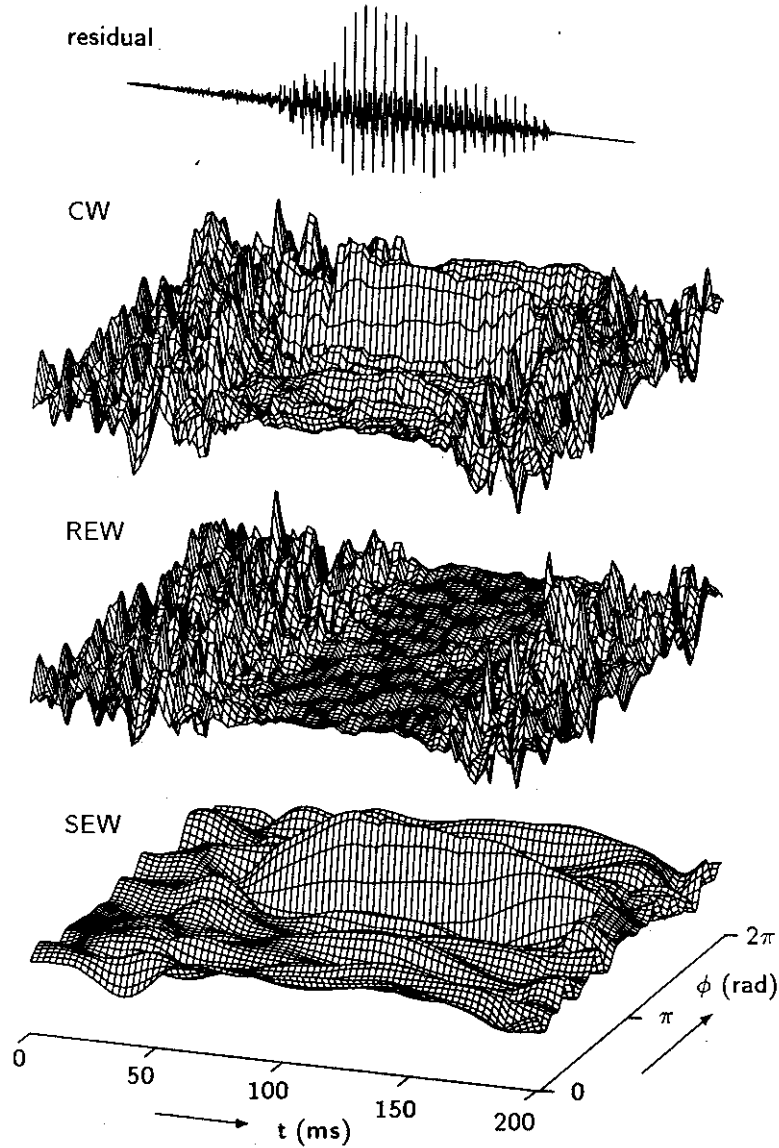
$$x_m[n] = u[n, t_m] \quad (7.84)$$

so that  $u[n, l]$  is defined for  $l = t_m$ , with the remaining points been computed through interpolation. A frequency representation of the pitch cycle can also be used instead of the time pitch cycle.

This surface can then be sampled at regular time intervals  $l = sT$ . It has been shown empirically that transmitting the pitch waveform  $x_s[n]$  about 40 times per second (a 25-ms interval is equivalent to  $T = 200$  samples for an  $F_s = 8000$  Hz sampling rate) is sufficient for voiced speech. The so-called *slowly evolving waveform* (SEW)  $\tilde{u}[n, l]$  can be generated by low-pass filtering  $u[n, l]$  along the  $l$ -axis:

$$x_s[n] = \tilde{u}[n, sT] = \frac{\sum_m h[sT - t_m] u[n, t_m]}{\sum_m h[sT - t_m]} \quad (7.85)$$

where  $h[n]$  is a low-pass filter and  $x_s[n]$  is a sampled version of  $\tilde{u}[n, l]$ .



**Figure 7.10** LP residual signal and its associated surface  $u(t, \phi)$ . In the  $\phi$  axis we have a normalized pitch pulse at every given time  $t$ . Decomposition of the surface into a slowly evolving waveform and a rapidly evolving waveform (After Kleijn [30], reprinted by permission of IEEE).



The decoder has to reconstruct each pitch waveform  $x_m[n]$  from the SEW  $x_s[n]$  by interpolation between adjacent pitch waveforms, and thus the name *waveform interpolation (WI) coding*:

$$\tilde{w}_m[n] = \tilde{u}[n, t_m] = \frac{\sum_s h[t_m - sT] w_s[n]}{\sum_s h[t_m - sT]} \quad (7.86)$$

If the sampling period is larger than the local pitch period ( $T > P_m$ ), perfect reconstruction will not be possible, and there will be some error in the approximation

$$x_m[n] = \tilde{x}_m[n] + \hat{x}_m[n] \quad (7.87)$$

or alternatively in the two-dimensional representation

$$u[n, l] = \tilde{u}[n, l] + \hat{u}[n, l] \quad (7.88)$$

where  $\hat{x}_m[n]$  and  $\hat{u}[n, l]$  represent the *rapidly evolving waveforms* (REW).

Since this technique can also be applied to unvoiced speech, where the concept of pitch waveform doesn't make sense, the more general term *characteristic waveform* is used instead. For unvoiced speech, an arbitrary *period* of around 100 Hz can be used.

For voiced speech, we expect the rapidly varying waveform  $\hat{u}[n, l]$  in Eq. (7.88) to have much less energy than the slowly evolving waveform  $\tilde{u}[n, l]$ . For unvoiced speech the converse is true:  $\hat{u}[n, l]$  has more energy than  $\tilde{u}[n, l]$ . For voiced fricatives, both components may be comparable and thus we want to transmit both.

In Eqs. (7.85) and (7.86) we need to average characteristic waveforms that have, in general, different lengths. To handle this, all characteristic waveforms are typically normalized in length prior to the averaging operation. This length normalization is done by padding with zeros  $x_m[n]$  to a certain length  $M$ , or truncating  $x_m[n]$  if  $P_m > M$ . Another possible normalization is done via linear resampling. This decomposition is shown in Figure 7.10.

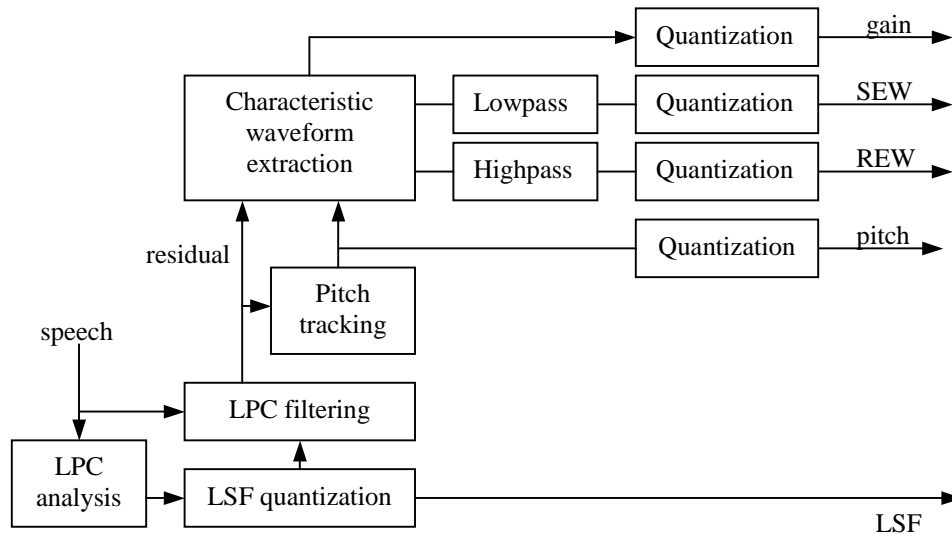
Another representation uses the Fourier transform of  $x_m[n]$ . This case is related to the harmonic model of Section 7.5.2. In the harmonic model, a relatively long window is needed to average the several pitch waveforms within the window, whereas this waveform interpolation method has higher time resolution. In constructing the characteristic waveforms we have implicitly used a rectangular window of length one pitch period, but other windows can be used, such as a Hanning window that covers two pitch periods. This frequency-domain representation offers advantages in coding both the SEW and the REW, because properties of the human auditory system can help reduce the bit rate. This decomposition is often done on the LPC residual signal.

In particular, the REW  $\hat{u}[n, l]$  has the characteristics for noise, and as such only a rough description of its power spectral density is needed. At the decoder, random noise is generated with the transmitted power spectrum. The spectrum of  $\hat{u}[n, l]$  can be vector quantized to as few as eight shapes with little or no degradation.

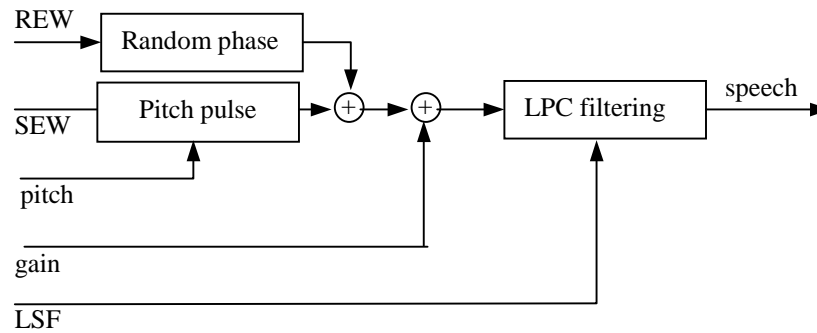
The SEW  $\tilde{u}[n, l]$  is more important perceptually, and for high quality the whole shape needs to be transmitted. Higher accuracy is desired at lower frequencies so that a perceptual

frequency scale (mel or Bark) is often used. Since the magnitude of  $\tilde{u}[n, l]$  is perceptually more important than the phase, for low bit rates the phase of the SEW is not transmitted. The magnitude spectrum can be quantized with the VDVQ described in Section 7.5.2.3.

To obtain the characteristic waveforms, the pitch needs to be computed. We can find the pitch period such that the energy of the REW is minimized. To do this we use the approaches described in Chapter 6. Figure 7.11 shows a block diagram of the encoder and Figure 7.12 of the decoder.



**Figure 7.11** Block diagram of the WI encoder.



**Figure 7.12** Block diagram of the WI decoder.

Parameter estimation using an analysis-by-synthesis framework [21] can yield better results than the open-loop estimation described above.

## 7.6. HISTORICAL PERSPECTIVE AND FURTHER READING

This chapter is only an introduction to speech and audio coding technologies. The reader is referred to [27, 32, 41, 52] for coverage in greater depth. A good source of the history of speech coding can be found in [20].

In 1939, Homer Dudley of AT&T Bell Labs first proposed the channel vocoder [15], the first analysis-by-synthesis system. This vocoder analyzed slowly varying parameters for both the excitation and the spectral envelope. Dudley thought of the advantages of bandwidth compression and information encryption long before the advent of digital communications.

PCM was first conceived in 1937 by Alex Reeves at the Paris Laboratories of AT&T, and it started to be deployed in the United States Public Switched Telephone Network in 1962. The digital compact disc, invented in the late 1960s by James T. Russell and introduced commercially in 1984, also uses PCM as coding standard. The use of  $\mu$ -law encoding was proposed by Smith [51] in 1957, but it wasn't standardized for telephone networks (G.711) until 1972. In 1952, Schouten et al. [47] proposed delta modulation and Cutler [11] invented differential PCM. ADPCM was developed by Barnwell [6] in 1974.

Speech coding underwent a fundamental change with the development of linear predictive coding in the early 1970s. Atal [3] proposed the LPC vocoder in 1971, and then CELP [5] in 1984. The majority of coding standards for speech signals today use a variation on CELP.

Sinusoidal coding [35] and waveform interpolation [29] were developed in 1986 and 1991, respectively, for low-bit-rate telephone speech. Transform coders such as MP3 [23], MPEG II, and Perceptual Audio Coder (PAC) [28] have been used primarily in audio coding for high-fidelity applications.

Recently, researchers have been improving the technology for cellular communications by trading off source coding and channel coding. For poor channels more bits are allocated to channel coding and fewer to source coding to reduce dropped calls. Scalable coders that have different layers with increased level of precision, or bandwidth, are also of great interest.

## REFERENCES

- [1] Adoul, J.P., *et al.*, "Fast CELP Coding Based on Algebraic Codes," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1987, Dallas, TX pp. 1957-1960.
- [2] Atal, B.S., R.V. Cox, and P. Kroon, "Spectral Quantization and Interpolation for CELP Coders," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1989, Glasgow pp. 69-72.
- [3] Atal, B.S. and L. Hanauer, "Speech Analysis and Synthesis by Linear Prediction of the Speech Wave," *Journal of the Acoustical Society of America*, 1971, **50**, pp. 637-655.
- [4] Atal, B.S. and M.R. Schroeder, "Predictive Coding of Speech Signals and Subjective Error Criteria," *IEEE Trans. on Acoustics, Speech and Signal Processing*, 1979, **ASSP-27**(3), pp. 247-254.
- [5] Atal, B.S. and M.R. Schroeder, "Stochastic Coding of Speech at Very Low Bit Rates," *Proc. Int. Conf. on Comm.*, 1984, Amsterdam pp. 1610-1613.

- [6] Barnwell, T.P., *et al.*, *Adaptive Differential PCM Speech Transmission*, 1974, Rome Air Development Center.
- [7] Benvenuto, N., G. Bertocci, and W.R. Daumer, "The 32-kbps ADPCM Coding Standard," *AT&T Technical Journal*, 1986, **65**, pp. 12-22.
- [8] Campbell, J.P., T.E. Tremain, and V.C. Welch, "The DoD 4.8 kbps Standard (Proposed Federal Standard 1016)" in *Advances in Speech Coding*, B. Atal, V. Cuperman, and A. Gersho, eds. 1991, pp. 121-133, Kluwer Academic Publishers.
- [9] Chen, J.H., *et al.*, "A Low-Delay CELP Coder for the CCITT 16 kbps Speech Coding Standard," *IEEE Journal on Selected Areas Communications*, 1992, **10**(5), pp. 830-849.
- [10] Chen, J.H. and A. Gersho, "Adaptive Postfiltering for Quality Enhancement of Coded Speech," *IEEE Trans. on Speech and Audio Processing*, 1995, **3**(1), pp. 59-71.
- [11] Cutler, C.C., *Differential Quantization for Communication Signals*, , 1952, US Patent 2,605,361.
- [12] Das, A. and A. Gersho, "Variable Dimension Vector Quantization," *IEEE Signal Processing Letters*, 1996, **3**(7), pp. 200-202.
- [13] Daumer, W.R., *et al.*, "Overview of the 32kbps ADPCM Algorithm," *Proc. IEEE Global Telecomm*, 1984 pp. 774-777.
- [14] DeJaco, P.J.A., W. Gardner, and C. Lee, "QCELP: The North American CDMA Digital Cellular Variable Speech Coding Standard," *Proc. Workshop on Speech Coding for Telecommunications*, 1993, Sainte Adele, Quebec pp. 5-6.
- [15] Dudley, H., "The Vocoder," *Bell Labs Record*, 1939, **17**, pp. 122-126.
- [16] Erdmann, C., *et al.*, "An Adaptive Rate Wideband Speech Codec with Adaptive Gain Re-Quantization," *IEEE Workshop on Speech Coding*, 2000, Delavan, Wisconsin.
- [17] Etemoglu, C.O., V. Cuperman, and A. Gersho, "Speech Coding with an Analysis-by-Synthesis Sinusoidal Model," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1371-1374.
- [18] Gerson, I.A. and M.A. Jasiuk, "Vector Sum Excited Linear Prediction (VSELP)" in *Advances in Speech Coding*, B.S. Atal, V. Cuperman, and A. Gersho, eds. 1991, Boston, MA, pp. 69-79, Kluwer Academic Publishers.
- [19] Gerson, I.A. and M.A. Jasiuk., "Techniques for Improving the Performance of CELP-type Speech Coders," *IEEE Journal Selected Areas Communications*, 1991, **10**(5), pp. 858-865.
- [20] Gold, B. and N. Morgan, *Speech and Audio Signal Processing: Processing and Perception of Speech and Music*, 2000, New York, John Wiley.
- [21] Gottesman, O. and A. Gersho, "High Quality Enhanced Waveform Interpolative Coding at 2.8 kbps," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1363-1366.
- [22] Greefkes, J.A., "A Digitally Companded Delta Modulation Modem for Speech Transmission," *Proc. Int. Conf. on Communications*, 1970 pp. 7.33-7.48.
- [23] ISO, *Coding of Moving Pictures and Associated Audio - Audio Part*, 1993, Int. Standards Organization.
- [24] ISO/IEC, *Information Technology - Coding of Moving Pictures and Associated Audio for Digital Storage Media at up to about 1.5 Mbps, Part 3: Audio (MPEG-1)*, 1992, Int. Standards Organization.
- [25] ITU-T, *Methods for Subjective Determination of Transmission Quality*, 1996, Int. Telecommunication Unit.
- [26] Jarvinen, K., *et al.*, "GSM Enhanced Full Rate Speech Codec," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1997, Munich, Germany pp. 771-774.
- [27] Jayant, N.S. and P. Noll, *Digital Coding of Waveforms*, 1984, Upper Saddle River, NJ, Prentice Hall.

- [28] Johnston, J.D., *et al.*, "ATT Perceptual Audio Coding (PAC)" in *Audio Engineering Society (AES) Collected Papers on Digital Audio Bit Rate Reduction*, N. Gilchrist and C. Grewin, eds. 1996, pp. 73-82.
- [29] Kleijn, W.B., "Continuous Representations in Linear Predictive Coding," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1991, Toronto, Canada pp. 201-204.
- [30] Kleijn, W.B. and J. Haagen, "Transformation and Decomposition of the Speech Signal for Coding," *IEEE Signal Processing Letters*, 1994, **1**, pp. 136-138.
- [31] Kleijn, W.B., D.J. Krasinski, and R.H. Ketchum, "An Efficient Stochastically Excited Linear Predictive Coding Algorithm for High Quality Low Bit Rate Transmission of Speech," *Speech Communication*, 1988, **7**, pp. 305-316.
- [32] Kleijn, W.B. and K.K. Paliwal, *Speech Coding and Synthesis*, 1995, Amsterdam, Netherlands, Elsevier.
- [33] Koishida, K., V. Cuperman, and A. Gersho, "A 16-KBIT/S Bandwidth Scalable Audio Coder Based on the G.729 Standard," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1149-1152.
- [34] Li, C. and V. Cuperman, "Analysis-by-Synthesis Multimode Harmonic Speech Coding at 4 kbps," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1367-1370.
- [35] McAulay, R.J. and T.F. Quateri, "Speech Analysis/Synthesis Based on a Sinusoidal Representation," *IEEE Trans. on Acoustics, Speech and Signal Processing*, 1986, **34**, pp. 744-754.
- [36] McAulay, R.J. and T.F. Quateri, "Sinusoidal Coding" in *Speech Coding and Synthesis*, W.B. Kleijn and K.K. Paliwal, eds. 1995, pp. 121-174, Elsevier.
- [37] McCree, A., "A 14 kbps Wideband Speech Coder with a Parametric Highband Model," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1153-1156.
- [38] McCree, A.V. and T.P. Barnwell, "Improving the Performance of a Mixed-Excitation LPC Vocoder in Acoustic Noise," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1992, San Francisco pp. II-137-138.
- [39] McCree, A.V., *et al.*, "A 2.4 kbit/s MELP Coder Candidate for the New U.S. Federal Standard," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1996, Atlanta, GA pp. 200-203.
- [40] Paez, M.D. and T.H. Glisson, "Minimum Squared-Error Quantization in Speech," *IEEE Trans. on Comm*, 1972, **20**, pp. 225-230.
- [41] Painter, T. and A. Spanias, "A Review of Algorithms for Perceptual Coding of Digital Audio Signals," *Proc. Int. Conf. on DSP*, 1997 pp. 179-205.
- [42] Painter, T. and A. Spanias, "Perceptual Coding of Digital Audio," *Proc. of IEEE*, 2000(April), pp. 451-513.
- [43] Paliwal, K.K. and B. Atal, "Efficient Vector Quantization of LPC Parameters at 24 Bits/Frame," *IEEE Trans. on Speech and Audio Processing*, 1993, **1**(1), pp. 3-14.
- [44] Prevez, M.A., H.V. Sorensen, and J.V.D. Spiegel, "An Overview of Sigma-Delta Converters," *IEEE Signal Processing Magazine*, 1996, **13**(1), pp. 61-84.
- [45] Robinson, T., *Simple Lossless and Near-Lossless Waveform Compression*, 1994, Cambridge University Engineering Department.
- [46] Salami, R., C. Laflamme, and B. Bessette, "Description of ITU-T Recommendation G.729 Annex A: Reduced Complexity 8 kbps CS-ACELP Codec," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1997, Munich, Germany pp. 775-778.
- [47] Schouten, J.S., F.E. DeJager, and J.A. Greefkes, *Delta Modulation, a New Modulation System for Telecommunications*, 1952, Phillips, pp. 237-245.

- [48] Shlomot, E., V. Cuperman, and A. Gersho, "Combined Harmonic and Waveform Coding of Speech at Low Bit Rates," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1998, Seattle, WA pp. 585-588.
- [49] Singhal, S. and B.S. Atal, "Improving Performance of Multi-Pulse LPC Coders at Low Bit Rates," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1984, San Diego pp. 1.3.1-1.3.4.
- [50] Skoglund, J., R. Cox, and J. Collura, "A Combined WI and MELP Coder at 5.2KBPS," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1387-1390.
- [51] Smith, B., "Instantaneous Companding of Quantized Signals," *Bell Systems Technical Journal*, 1957, **36**(3), pp. 653-709.
- [52] Spanias, A.S., "Speech Coding: A Tutorial Review," *Proc. of the IEEE*, 1994, **82**(10), pp. 1441-1582.
- [53] Stachurski, J. and A. McCree, "A 4 kbps Hybrid MELP/CELP Coder with Alignment Phase Encoding and Zero Phase Equalization," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1379-1382.
- [54] Todd, C., "AC-3: Flexible Perceptual Coding for Audio Transmission and Storage," *Audio Engineering Society 96th Convention*, 1994.
- [55] Tremain, T.E., *The Government Standard Linear Predictive Coding Algorithm*, in *Speech Technology Magazine*, 1982. pp. 40-49.
- [56] Vary, P., *et al.*, "A Regular-Pulse Excited Linear Predictive Code," *Speech Communication*, 1988, **7**(2), pp. 209-215.
- [57] Wang, T., *et al.*, "A 1200 BPS Speech Coder Based on MELP," *Int. Conf. on Acoustics, Speech and Signal Processing*, 2000, Istanbul, Turkey pp. 1375-1378.